

STEX+ – a System for Flexible Formalization of Linked Data

Andrea Kohlhase
German Research Center for
Artificial Intelligence (DFKI)
Enrique-Schmidt-Str. 5
28359 Bremen, Germany
Andrea.Kohlhase@dfki.de

Michael Kohlhase
Jacobs University Bremen
P. O. Box 750561
28725 Bremen, Germany
m.kohlhase@jacobs-
university.de

Christoph Lange
Jacobs University Bremen
P. O. Box 750561
28725 Bremen, Germany
ch.lange@jacobs-
university.de

ABSTRACT

We present the STEX+ system, a user-driven advancement of STEX — a semantic extension of L^AT_EX that allows for producing high-quality PDF documents for (proof)reading and printing, as well as semantic XML/OMDoc documents for the Web or further processing. Originally STEX had been created as an invasive, semantic frontend for authoring XML documents. Here, we used STEX in a Software Engineering case study as a formalization tool. In order to deal with modular pre-semantic vocabularies and relations, we upgraded it to STEX+ in a participatory design process. We present a tool chain that starts with an STEX+ editor and ultimately serves the generated documents as XHTML+RDFa Linked Data via an OMDoc-enabled, versioned XML database. In the final output, all structural annotations are preserved in order to enable semantic information retrieval services.

Categories and Subject Descriptors

D.2.1 [Software Engineering]: Requirements/Specifications—*Languages*; I.2.4 [Artificial Intelligence]: Knowledge Representation Formalisms and Methods—*Representation languages*; I.7.2 [Document and Text Processing]: Document Preparation

General Terms

Documentation, Human Factors, Languages, Management

Keywords

formalization, L^AT_EX, Linked Data, software engineering, semantic authoring, annotation, metadata, RDFa, vocabularies, ontologies

1. INTRODUCTION

An important issue in the Semantic Web community was and still is the “Authoring Problem”: How can we convince people not only to use semantic technologies, but also prepare them for creating semantic documents (in a broad sense)?

Here, we were interested in formalizing a collection of L^AT_EX documents into a set of files in the OMDoc format, an XML vocabulary specialized for managing mathematical information, and further on to Linked Data for interactive browsing and querying on the Semantic Web.

Concretely, the object of our study was the collection of documents created in the course of the 3-year project “Sicherungskomponente für Autonome Mobile Systeme (SAMS)” at the German Research Center for Artificial Intelligence (DFKI). SAMS built a software safety component for autonomous mobile service robots developed and certified it as SIL-3 standard compliant (see [13]). Certification required the software development to follow the V-model (figure 1) and to be based on a verification of certain safety properties in the proof checker Isabelle [33]. The V-model mandates e.g. that relevant document fragments get justified and linked to corresponding fragments in other members of the document collection in an iterative refinement process (the arms of the ‘V’ from the upper left over the bottom to the upper right and in-between in figure 1).

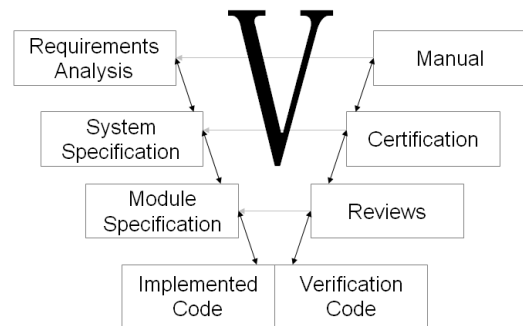


Figure 1: A Document View on the V-Model

System development with respect to this regime results in a highly interconnected collection of design documents, certification documents, code, formal specifications, and formal proofs. This collection of documents “SAMSdocs” [35] make up the basis of a case study in the context of the FormalSafe project [12] at DFKI Bremen, where they serve as a basis for research on machine-supported change management, information retrieval, and document interaction. In this paper, we report on the formalization project of the collection of L^AT_EX documents in SAMSdocs (that we will without further ado also abbreviate with SAMSdocs).

Not surprisingly, the interplay between the fields Semantic Web and Human-Computer Interaction played an important role as the “Authoring Problem” of the first is often tackled via methods of the second. One such approach is that of “*invasive technology*” [21] with the basic idea that from a user’s perspective, semantic authoring and general editing are the same, so why not offer semantic functionalities as an extension of well-known editing systems, thereby ‘invading’ the existent ones. We started with L^AT_EX not only because a good portion of our case study was written in it, but also as L^AT_EX constitutes the state-of-the-art authoring solution for many scientific/technical/mathematical document collections. Despite its text-based nature it is widely considered the most efficient tool for the task. Therefore, we used the invasive OMDoc frontend for L^AT_EX documents called **gT_EX** [26]. In the formalization process its conceptual usability weaknesses (for the task) were identified and within a participatory design process it evolved into the invasive formalization tool **gT_EX+**.

In section 2, we will present the **gT_EX** system, especially its realization of Linked Data creation. Then we describe in section 3 the formalization process of SAMSDocs with **gT_EX**, our challenges, and our (pre-)solutions. In section 4 we report the enhancements of **gT_EX** realized in and for the case study to **gT_EX+**. Having **gT_EX+** documents with Linked Data and ontological markup, we describe (potential) services and their implementation design in section 5. Section 6 summarizes related work, and section 7 concludes the paper.

2. **gT_EX**: OBJ.-ORIENTED L^AT_EX MARKUP

gT_EX [26, 37] is an extension of the L^AT_EX language that is geared towards marking up the semantic structure underlying a document. The main concept in **gT_EX** is that of a “**semantic macro**”, i.e., a T_EX command sequence \mathcal{S} that represents a meaningful (mathematical) concept \mathcal{C} : the T_EX formatter will expand \mathcal{S} to the presentation of \mathcal{C} . For instance, the command sequence `\positiveReals` (from listing 1) is a semantic macro that represents a mathematical symbol — the set \mathbb{R}^+ of positive real numbers. While the use of semantic macros is generally considered a good markup practice for scientific documents (e.g., because they allow to adapt notation by macro redefinition and thus increase reusability), regular T_EX/L^AT_EX does not offer any infrastructural support for this. **gT_EX** does just this by adopting a semantic, ‘object-oriented’ approach to semantic macros by grouping them into “modules”, which are linked by an “imports” relation. To get a better intuition, consider

Listing 1: An **gT_EX Module for Real Numbers**

```
\begin{module}[id=reals]
  \importmodule[../background/sets]{sets}
  \symdef{Reals}{\mathbb{R}}
  \symdef{greater}[2]{\#1>\#2}
  \symdef{positiveReals}{\Reals^+}
  \begin{definition}[id=posreals.def,
    title=Positive Real Numbers]
    \defeq\positiveReals
      {\setst{\inset{x}\Reals}{\greater{x}0}}$
  \end{definition}
  ...
\end{module}
```

which would be formatted to

Definition 2.1 (Positive Real Numbers): $\mathbb{R}^+ := \{x \in \mathbb{R} \mid x > 0\}$

Here, **gT_EX**’s `\symdef` macro – invasive by to its deliberate resemblance of (L^A)T_EX’s `\def` and `\newcommand` – generates a respective semantic macro, for instance the `\positiveReals` with representation \mathbb{R}^+ . Note the symbol inheritance scheme of **gT_EX**: The markup in the module `reals` has access to semantic macros `\setst` (“set such that”) and `\inset` (set membership) from the module `sets` that was imported by the document `\importmodule` directive from the `../background/sets.tex`. Furthermore, it has access to the `\defeq` (definitional equality) that was in turn imported by the module `sets`.

From this example we can already see an organizational advantage of **gT_EX** over L^AT_EX: we can define the (semantic) macros close to where the corresponding concepts are defined, and we can (recursively) import mathematical modules. But the main advantage of markup in **gT_EX** is that it can be transformed to XML via the L^AT_EXML system [32]: Listing 2 shows the OMDoc [25] representation generated from the **gT_EX** sources in listing 1. **OMDoc** is a semantics-oriented representation format for mathematical knowledge that extends the formula markup formats OpenMath [7] and MathML [2] to a document markup format.

Listing 2: An XML Version of Listing 1

```
<theory xml:id="reals">
  <imports from="../background/sets.omdoc#sets"/>
  <symbol xml:id="Reals"/>
  <notation>
5    <prototype><OMS cd="reals" name="Reals"/></prototype>
    <rendering><m:mo>\mathbb{R}</m:mo></rendering>
  </notation>
  <symbol xml:id="greater"/><notation>...</notation>
  <symbol xml:id="positiveReals"/><notation>...</notation>
10  <definition xml:id="posreals.def" for="positiveReals">
    <meta property="dc:title">Positive Real Numbers</meta>
    <OMOBJ>
      <OMA>
        <OMS cd="mathtalk" name="defeq"/>
        <OMS cd="reals" name="positiveReals"/>
15      <OMA>
        <OMS cd="sets" name="setst"/>
        <OMA>
          <OMS cd="sets" name="inset"/>
          <OMV name="x"/>
20          <OMS cd="reals" name="reals"/>
        </OMA>
        <OMA>
          <OMS cd="reals" name="greater"/>
          <OMV name="x"/>
25          <OMI>0</OMI>
        </OMA>
      </OMA>
    </OMOBJ>
  </definition>
  ...
</theory>
```

One thing that jumps out from the XML in this listing is that it incorporates all the information from the **gT_EX** markup that was invisible in the PDF produced by formatting it with T_EX.

OMDoc itself has been used as a storage and exchange format for automated theorem provers, software verification systems, e-learning software, and other applications [25, chap-

ter 26], but due to its focus on semantic structures, it is not intended to be consumed by human readers. The Java-based JOMDoc [19] library uses the `notation` elements to generate human-readable XHTML+MathML from OMDoc. Figure 2 shows the result of rendering the document from listing 2 in a MathML-aware browser. In contrast to the PDF output we can directly create from \LaTeX , XHTML+MathML allows for interactivity. In particular, our JOBAD JavaScript framework enables modular interactive services in rendered XHTML+MathML documents [14]. These services utilize the semantic structures of mathematical formulae. In our rendered documents, each formula in human-readable Presentation MathML carries the original semantic OpenMath representation of the formula, as shown in listing 2, as a hidden annotation.

Client-side JOBAD services, which exclusively rely on annotations given inside a document, have already been implemented for folding and unfolding subterms of formulae and for controlling the display of redundant brackets in complex formulae. The symbol definition lookup service, shown in figure 2, interacts with a server backend: It traverses the links to symbol and their corresponding definition elements that are established by the OMS elements in OpenMath – for example, `<OMS cd="sets" name="inset"/>` encodes the URI `../background/sets.omdoc#inset` – and retrieves the document at that URI as XHTML+MathML.¹ JOBAD’s ability to integrate an arbitrary number of services, which can talk to different server backends and which are enabled depending on the context, i.e., the semantic structure of the part of a mathematical formula that the user has selected, turns our rendered mathematical documents into powerful mashups [28]. On any symbol, for example, definition lookup is enabled. On any expression where a number is multiplied with a special symbol representing a unit of measurement, a unit conversion client that talks to a remote unit conversion web service is enabled. The JOBAD architecture has been designed without depending on a particular backend; for most of our services we are using the extensible XML-aware database TNTBase [39, 40, 11], which has special support for OMDoc and integrates the JOMDoc rendering library.

DEFINITION:

$\mathbb{R}^+ := \{x \in \mathbb{R} \mid x > 0\}$

Definition Lookup (defeq)

The simplest form of definition schema is the **simple definition**. This just introduces a name (the **definiendum**) for a compound object (the **definiens**). Note that the name must be new, i.e. may not have been used for anything else, in particular, the definiendum may not occur in the definiens. We use the symbols `:=` (and the inverse `=:`) to denote simple definitions in formulae.

Figure 2: Listing 1 as Dynamic XHTML+MathML

¹This is the MathML way of representing Linked Data. In section 5, we describe how we have now extended this feature to cover RDFa Linked Data.

3. FORMALIZATION WITH \LaTeX TOWARDS $\text{\LaTeX}+$

In this section we describe the process of formalizing the SAMSDocs collection of \LaTeX documents created in the course of the SAMS project with the \LaTeX system. We use the user’s perspective to point to the requirements for $\text{\LaTeX}+$ that evolved in this process.

As we all know all too well: Formalizing is never easily done. In our project we had the additional challenge of doing it without corruption of the PDF layout that was produced with \LaTeX . Here, \LaTeX fits well, as it generates PDF and transforms to XML. In figure 3 we can see the general course of action:

- i) we identified document fragments (“**objects**”) that constitute a coherent, meaningful unit like the state of a document “rd.” or its description “ready for certification”, then
- ii) we translated it into the \LaTeX format, realizing for example that “rd.” is a recurring symbol and “ready for certification” its definition (therefore designing the SAMSDocs macro “SDdef”), and finally
- iii) we polished these macros in the \LaTeX specific sty-files so that the PDF layout remained as before and the generated XML represented the intended logical structure, for instance the use of the OMDoc XML elements symbol and definition.

Note that definitions are common objects in mathematical documents, therefore \LaTeX naturally provides a definition environment. So why didn’t we use that? Because the document model of OMDoc, which we obtain by transforming \LaTeX using \LaTeX XML, does not allow definitions in tables, as the former are stand-alone objects from an ontological perspective. If one *authors* a formal document, this view is taken, so no problem arises, but if one *formalizes* an existing document, layout and cognitive side-conditions have to be taken into account. We therefore realized that we could not simply add basic \LaTeX markup to the \LaTeX source yielding formal objects, we rather needed to add pre-formal markup in the formalization process (we speak of (**semantic**) **preloading**).

Whenever project-wide (semantic) layout schemes were discovered, that were frequently used, we extended the macro set of \LaTeX suitably (enabling preloading “*project structures*” [22], i.e. project-induced ones which is quite different from “*document [layout] structures*” [ibid.], e.g. by subsections that is supported by \LaTeX core features, see DCMsubsection in figure 3). The table layout for example was often used for lists of symbol definitions. So we created the SDTab-def environment which can host as many SDdef commands as wanted (see fig. 3). This increased the efficiency of the formalizing process tremendously.

Another difference between authoring and semantic preloading consisted in the *order of the formalization steps*. While the order of the first typically consists of “**chunking**” (i.e., building up structure e.g. by setting up theories), “**spotting**” (i.e., coining objects), and “**relating**” (i.e., making relationships between objects or structures explicit), the order of the second is made up of spotting, then relating *or* chunk-

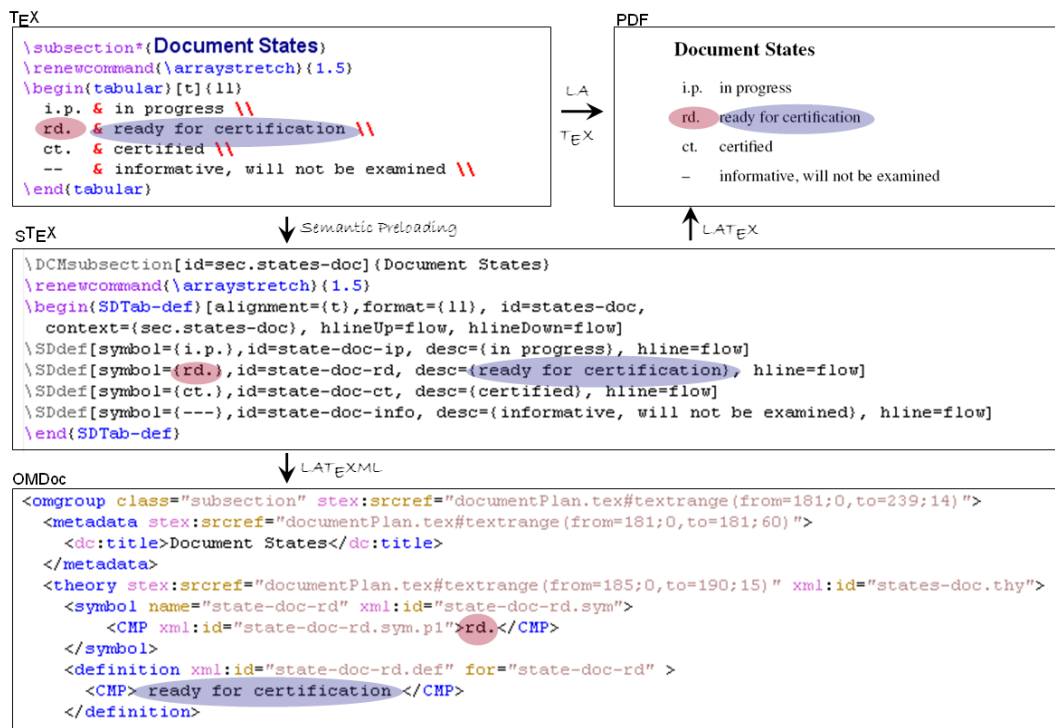


Figure 3: The Formalization Workflow via sTeX: Definition Table of “document state”

ing. The last two were done simultaneously, because sTeX offers a very handy inheritance scheme for symbol macros — as long as the chunks are in order, which could be sensibly done for some but not for all at this stage in the formalization process. Generally, many ‘guiding’ services of sTeX, that sTeX considered to be features, turned out to be too rigid.

As a consequence we heavily used very light annotations at the beginning: It was sufficient to identify a certain document fragment and to mark it with a referencable ID like “state-doc-rd”. Shortly afterwards, we realized that some more basic markup was necessary, since we wanted to formalize our knowledge of types/categories of these objects and their conceptual belonging. For this we developed a set of “ad-hoc semantification macros” with named attributes like SDObject[id], SDmore[id,cat,for], SDisa[id,cat,for,for, follows,theory,imports,tab], or SDreferences[id,file,refid]². The ‘more’ functionality provided by SDmore was required due to logically contiguous objects that were interspersed in a document. With this set we preloaded “object structures” [ibid.], i.e. object-induced ones. Note that the ad-hoc semantification macros enabled the formalizer to develop her own metadata vocabulary.

As soon as the document boundaries went down, we realized that an object had many occurrences in several of the documents in the SAMSDocs collection. For example, first

²We use subsets of a general attributes set for all of our sTeX extensions to lower the learning curve for the use of the markup macros.

an object was introduced as a high-level concept in the contract, then it was specified in another document, refined in a detailed specification, implemented in the code, reviewed at some stage, and so on until it was finally described in the manual. Thus, we had to preload “collection structures” [ibid.] as well, which consisted in the development process model, the V-model as seen in figure 1. Here, we built our personal **V-model macros**, e.g. SemVMrefines, SemVMimplements, or SemVMdescribesUse.

Additionally, we created an sTeX extension especially suited for preloading “organizational structures” [ibid.]. This is considered different from project structures as organizational markup is very probable to be reusable for other projects with the same organizational structures. For example, SAMS used a document version management as well as a document review history, so that environments VMchangelist, VMcertification with corresponding list entry macros VMchange, VMcertified were built. Another example is the processing state of a document, which can be marked up easily by using the VMdocstate macro as seen in figure 4.

We noted that the necessary formalization depth of some documents was naturally deeper than others. For example, it didn’t seem sensible to formalize the contract too much, as it was created as a high-level communication document, whereas the detailed specification needed a lot of formalization. The manual had an interesting mixed state of formality and informality, as it was again geared towards communication, but it needed to be very precise. In conclusion we note that the mathematical content of the documents (i.e., the mathematical objects and their relations) was only one of the knowledge sources that needed to be formalized and

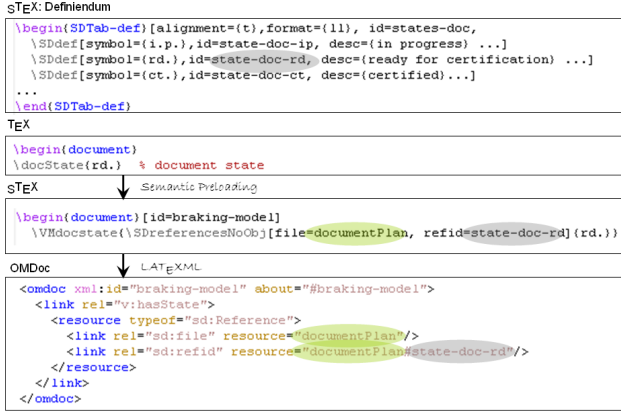


Figure 4: Referencing a “document state”

marked up. In the course of the formalization it has become apparent that the knowledge in such complex collections is *multi-dimensional* (cf. [22] for an in-depth analysis). Thus, the requirements for extending $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ to $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}+$ were (i) to generate XML output that preserves the semantics annotated in the preloading phase, (ii) and to take into account the multi-dimensionality of our ad-hoc semantification macros in a way that technically enables browsing and querying. These requirements were satisfied by enabling the generation of RDFa from our annotations and making them accessible to Linked Data services, as we will describe in the following sections.

4. $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}+$: A METADATA-EXTENSION OF $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$

All the arrows in figure 1 are examples of relations between document fragments in the $\mathcal{S}\mathcal{A}\mathcal{M}\mathcal{S}\mathcal{D}\mathcal{O}\mathcal{C}\mathcal{S}$ corpus that needed to be made explicit in addition to the mathematical relations that $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ had originally supported; the revision histories of documents and the social networks of their authors constitute further dimensions of knowledge. For situations like these, we had incorporated RDFa [1] as a flexible metadata framework into the OMDoc format [31]. In the course of this case study, the RDFa integration was revised and extended and will become part of the upcoming OMDoc version 1.3 [27]. The main idea for this integration is to realize that any concrete document markup format can only treat a certain set of objects and their relations via its respective native markup infrastructure. All other objects and relations can be added via RDFa annotations to the host language – assuming the latter is XML-based.

It is crucial to realize that, for machine support, the metadata objects and relations are given a machine-processable meaning via suitable ontologies. Moreover, ontologies are just special cases of (mathematical) theories, which import appropriate theories for the logical background, e.g. description logic, and whose symbols are the entities (class, properties, individuals) of ontologies. Thus, $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ and OMDoc can play a dual role for Linked Data in documents with mathematical content. They can be used as markup formats for the documents and at the same time as the markup formats for the ontologies. We have explored this correspondence

for OMDoc in previous work and implemented a translation between OMDoc and OWL [31, 30].

To understand our contribution, note that we can view $\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}$ and $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ as frameworks for defining domain-specific vocabularies in classes and packages; $\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}$ is used for layout aspects, and $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ can additionally handle the semantic aspects of the vocabularies. $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ uses this approach to define special markup e.g. for definitions (see lines 10 to 31 in listing 2). Note that to define $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ markup functionality like the definition environment, we have to provide a $\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}$ environment definition (so that the formatting via $\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}$ works) and a $\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}\mathcal{M}\mathcal{L}$ binding (to specify the XML transformation for the definition environment). As the OMDoc vocabulary is finite and fixed, $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ can (and does) supply special $\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}$ macros and environments and their $\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}\mathcal{M}\mathcal{L}$ bindings. But the situation is different for the flexible, RDFa-based metadata extension in OMDoc 1.3 we mentioned above, with a potentially infinite supply of vocabularies. At the start of the $\mathcal{S}\mathcal{A}\mathcal{M}\mathcal{S}\mathcal{D}\mathcal{O}\mathcal{C}\mathcal{S}$ preloading effort, $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ already supported a common subset of metadata vocabularies. For instance the Dublin Core title metadata element in line 11 of listing 2 is the transformation result of using the KeyVal [9] pair title=... in the optional argument of the definition environment.

For the $\mathcal{S}\mathcal{A}\mathcal{M}\mathcal{S}\mathcal{D}\mathcal{O}\mathcal{C}\mathcal{S}$ case study we started in the same way by adding a package with $\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}\mathcal{M}\mathcal{L}$ bindings to $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$. The $\backslash\mathcal{V}\mathcal{M}\mathcal{D}\mathcal{O}\mathcal{C}\mathcal{S}\mathcal{T}\mathcal{A}\mathcal{T}\mathcal{E}$ macro shown in the “ $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ ” box of figure 4 allowed us to annotate a document with its processing state. This is transformed to an RDFa-annotated omdoc root element, as shown in the “OMDoc” box underneath and in the black, solid parts of the RDF graph in figure 5. We can already see that the $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ extension for $\mathcal{S}\mathcal{A}\mathcal{M}\mathcal{S}\mathcal{D}\mathcal{O}\mathcal{C}\mathcal{S}$ exactly consists in a domain-specific metadata vocabulary extension, and that using the custom vocabulary hides markup complexity from the author. Again, $\mathcal{S}\mathcal{A}\mathcal{M}\mathcal{S}\mathcal{D}\mathcal{O}\mathcal{C}\mathcal{S}$ only needed a finite vocabulary extension, so this approach was feasible, but of restricted applicability, since developing the $\mathcal{S}\mathcal{A}\mathcal{M}\mathcal{S}\mathcal{D}\mathcal{O}\mathcal{C}\mathcal{S}$ package for $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ required insights into $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ internals and $\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}\mathcal{M}\mathcal{L}$ bindings. Thus this extension approach lacks the flexible user-extensibility that would be needed to scale up further.

To enable user-extensibility, we add a new declaration form $\backslash\mathcal{K}\mathcal{E}\mathcal{Y}\mathcal{D}\mathcal{E}\mathcal{F}$ to the core $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ functionality (yielding $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}+$) — like $\backslash\mathcal{S}\mathcal{Y}\mathcal{M}\mathcal{D}\mathcal{E}\mathcal{F}$ in that it is inherited via the module imports relation, only that it defines a KeyVal key instead of a semantic macro. To understand its application, we rationally reconstruct the $v:hasState$ relation from the example in the OMDoc box of figure 4. To do this, we use $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ to create a metadata vocabulary for document states: we create a certification module, which defines the $hasState$ metadata relation and adds it to the KeyVal keys of the document environment. The metalanguage macro is a variant of $\mathcal{I}\mathcal{M}\mathcal{P}\mathcal{O}\mathcal{R}\mathcal{T}\mathcal{M}\mathcal{O}\mathcal{D}\mathcal{U}\mathcal{L}\mathcal{E}$ that imports the meta language, i.e., the language in which the meaning of the new symbols is expressed; here we use OWL.

Listing 3: A Metadata Ontology for Certification

```
\begin{module}[id=certification]
\metalinguage[../background/owl]{owl}
\keydef{document}{hasState}
```



```

\symdef{state-doc-rd}[1]{rd. #1}
5 \symdef{tuev}{\text{T"UV}}
\begin{definition}[for=hasState]
  A document {\definiendum[hasState]{has state}} $x$, iff
  the project manager decrees it so.
\end{definition}
10 \begin{definition}[for=state-doc-rd]
  A document has state \definiendum[state-doc-rd]{rd. $x$},
  iff it has been submitted to $x$ for certification.
\end{definition}
\begin{definition}[for=tuev,hasState=$\statedocrd\tuev$]
15 The $\tuev$ (Technischer \Überwachungsverein) is a
  well-known certification agency in Germany.
\end{definition}
\end{module}

```

In this paper, we focus on using $\text{\texttt{\textit{S}T\textit{E}X+}}$ as a language for defining lightweight vocabularies. Note, however, that “heavyweight” formal semantics can be added to vocabulary terms in the same way as has been shown for mathematical symbols in listing 1. Similarly as the “real numbers” module relies on an $\text{\texttt{\textit{S}T\textit{E}X}}$ module that introduces set theory, the certification ontology relies on an $\text{\texttt{\textit{S}T\textit{E}X}}$ module that introduces the OWL language. Such an OWL ontology that has been written in $\text{\texttt{\textit{S}T\textit{E}X+}}$ can be translated to one of the widely supported serializations of OWL via two paths: (i) In the original workflow, the $\text{\texttt{\textit{S}T\textit{E}X+}}$ source is translated to OMDoc. Thanks to their modularity and literal programming capabilities, the $\text{\texttt{\textit{S}T\textit{E}X+}}$ or OMDoc representation allows for an expressive documentation of OWL ontologies. But, as OMDoc is not universally understood on the Semantic Web, we have implemented a translation of OWL ontologies encoded and documented in OMDoc to the standard RDF/XML representation [31]. (ii) Alternatively to this previously existing translation via OMDoc as an intermediate representation, we are working on a direct $\text{\texttt{\textit{S}T\textit{E}X+}}$ to OWL transformation. Simply using our experimental `owl2onto` class [23] instead of the `omdoc` class from $\text{\texttt{\textit{S}T\textit{E}X}}$ in the $\text{\texttt{\textit{L}A\textit{T}\textit{E}X}}$ preamble will cause $\text{\texttt{\textit{L}A\textit{T}\textit{E}X}}$ to generate OWL – here in the direct OWL XML serialization – instead of OMDoc from a subset of the $\text{\texttt{\textit{S}T\textit{E}X+}}$ markup.

Listing 4: Annotating a Document with Certification Metadata

```

\importmodule{../ontologies/cert}{certification}
2 \begin{document}[hasState=$\statedocrd\tuev$]
...
\end{document}

```

Let us now see how to *use* a vocabulary: If we import the certification metadata module, we can write to generate RDFa annotations that correspond to the (red) dotted arrow in figure 5. Note that in the state of formalization shown in figure 4, the SAMSDocs-specific RDF vocabulary still has a pre-semantic structure. With the $\text{\texttt{\textit{S}T\textit{E}X+}}$ we can express that the processing state is actually intended to be a symbol in a metadata theory, not just some semantic object in some file. In listing 3 we use the `\symdef` directive to generate the symbol `state-doc-rd` and `\keydef` to generate a metadata relation `hasState` that is expressed by a key of the same name, which is added to the document environment. When processed by $\text{\texttt{\textit{L}A\textit{T}\textit{E}X}}$, `\keydef` takes care of generating correct URIs for the metadata relations and their target resources, resulting in an RDFa output syntactically similar to figure 4. In conclusion, we note that

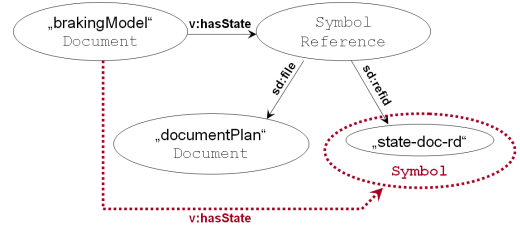


Figure 5: RDF View on a “doc. state” Assignment

$\text{\texttt{\textit{S}T\textit{E}X+}}$ allows us to rationally recreate the effect we previously achieved with the custom `\VMdocstate` and `\SD referencesNoObj` macros. Note that we did not have to extend the $\text{\texttt{\textit{L}A\textit{T}\textit{E}X}}$ bindings at all for this extension. Thus, $\text{\texttt{\textit{S}T\textit{E}X+}}$ gives us a generic $\text{\texttt{\textit{T}\textit{E}X}} \rightarrow \text{RDFa}$ translation, which works for arbitrary vocabularies.³

5. $\text{\texttt{\textit{S}T\textit{E}X+}}$ DOCUMENTS AS LINKED DATA

The translation of classical $\text{\texttt{\textit{S}T\textit{E}X}}$ to OMDoc and further to XHTML+MathML (see section 2), which results in a Linked Data like markup for mathematical symbols, enables interactive services in mathematical formulae. Now that $\text{\texttt{\textit{S}T\textit{E}X+}}$ supports formalization with arbitrary metadata (cf. section 4), it should additionally be possible to utilize *these* metadata for services. Both types of annotation complement each other: A practical $\text{\texttt{\textit{S}T\textit{E}X+}}$ document, like many of the SAMSDocs, would combine elements from listing 4 with those from listing 1 and consequently rely on services for both types of semantic structures.

The JOBAD service architecture (see section 2) gives uniform access to common queries in the document browsing user interface. In the SAMSDocs scenario this might be a query for all persons who have worked on the current document. This can directly be answered from the metadata of the revision log. Another typical query would consist in asking for all parts of a specification that have to be re-certified. Answering this query involves revision logs (for finding documents that have changed since the last certification), collection structures (V-model dependencies of changed parts), and mathematical structures (logical dependencies). In [22] we have elaborated on such SAMSDocs queries from the point of view of their stakeholders (like engineers, project managers, certifiers), particularly exploring the multi-dimensionality of the formal structures. For example, a project manager may find a substitute for an employee *E*, who has implemented a specification, by tracing back a link from the documentation of the implementation to the specification document and finding out, from the metadata of that document, who has recently been working on it. Here, we will summarize the extensions made to our system architecture to enable these services.

As a first step, we made the JOMDoc renderer preserve the RDFa metadata from the OMDoc documents, now gener-

³Our experimental `rdfameta` package [24] extends this to arbitrary $\text{\texttt{\textit{L}A\textit{T}\textit{E}X}}$ documents: It redefines common $\text{\texttt{\textit{L}A\textit{T}\textit{E}X}}$ commands (e.g. the sectioning macros) so that they include optional `KeyVal` arguments that can be extended by `\keydef` commands. With this metadata extension, we can add RDFa metadata to *any* existing $\text{\texttt{\textit{L}A\textit{T}\textit{E}X}}$.

ating XHTML+MathML+RDFa. Additionally, the mathematical structures (those that are above the formula level) had to be preserved in the rendered output. Even though OMDoc uses native non-RDFa markup for these structures, we can also represent these in RDF, exploiting the OMDoc ontology (see [29, 11] for more information). Existing JOBAD services recognized mathematical formulae in XHTML presentations of OMDoc documents by their semantic structure (e.g. whether they use previously defined symbols or units of measurement). Similarly, new services can now recognize from the RDFa annotations whether a chunk of an XHTML document is, e.g., an implementation of a specification fragment, and by which user requirement that is induced. Compared to the previously existing definition lookup service, the principle of retrieving content from a target URI and displaying it in a popup remained the same – the URIs are just provided by different annotations.

Secondly, we have extended the folding of subterms of mathematical formulae to higher-level structures, such as requirements or steps of structured proofs. We have implemented this using the `rdfQuery` JavaScript library [38], which parses all RDFa annotations of a document into a local triple store that can be queried using SPARQL-like JavaScript functions. On the server side, we have extended TNTBase [39], our versioned database backend and web server/application framework to accept commits of \LaTeX documents, automatically convert them to OMDoc, and then serve OMDoc, XHTML+MathML+RDFa, and, optionally, RDF/XML, according to the Linked Data best practices [17].

Even the pre-semantic annotations like the ones shown in figure 4 afford interactive services: A generic reference can already be utilized for lookup and navigation. Providing additional information in the instance document or in the ontology (e.g. the knowledge about the target of a reference being a symbol or a processing state) allows for making the service user interface more specific and enables the display of more relevant related information. For the generic pre-semantic “references” relation, the list of all semantic objects that it relates to each other would be too large for being usable, as there is no obvious way of ranking or filtering the link targets. But once more specific link types are used, such as the “has state” link, that information can be used to display a list of documents grouped by state.

Queries across documents cannot be answered using the above-mentioned `rdfQuery`: client side queries require a combination of querying a local triple store and crawling links. In our setup, we have experimented with SQUIN [16], a frontend to the Semantic Web Client library [4], which gives access to Linked Data via a simple HTTP frontend at very low integration costs: If the server provides standard-compliant Linked Data, then the client simply has to access the URL of the SQUIN server, providing a SPARQL query as a parameter. An alternative would have been AJAR library, a part of the Tabulator Linked Data browser [3], which implements the same functionality in JavaScript. In our test setup, SQUIN acted as a proxy between the client-side JavaScript code and our Linked Data. While a Linked Data crawler is most flexible when data are distributed across many servers (e.g. when an OMDoc document links to DBpedia), its query answering capabilities are only as good as

the Linked Data being served. For example, if the RDF(a) does not contain back-links (like links from a mathematical theory to the theories it imports *and* to the theories by which it is imported), then an AJAR- or SQUIN-powered client cannot query links in both directions. Moreover, the performance of such a solution is limited, as it requires memory for the local triple store as well processor time for query answering on the client side. Therefore, in the `SAMSDocs` setting, where the queries are currently limited to a document collection on a single server, the best solution is storing the triples on that same server, and making them accessible via a standard query interface. Concretely, we make a SPARQL endpoint powered by the Virtuoso triple store [34] available as an extension to TNTBase [11]. In a larger Software Engineering scenario (like a document collection of a company with multiple departments) a combination with a Linked Data crawler, as offered by the Sponger extension to Virtuoso in an integrated server-side fashion, may have advantages: if all these departments publish their document collections as Linked Data in the company intranet (see for instance [36] for the topicality of this example), crawling these may reveal previously unknown connections, e.g. colleagues dealing with structurally similar problems who could lend advice. Note that local vocabularies resulting from ad-hoc semantification need not be a barrier to knowledge exchange: Linked Data practices recommend connecting occurrences of semantically equivalent resources in different data sets by *owl:sameAs*. Alternatively, if it turns out that one department uses a “better” vocabulary for their data, the \LaTeX metadata extensions make it easy to adopt it: all we have to do is to change the \LaTeX bindings or `\keydefs`.⁴

6. RELATED WORK

We have presented \LaTeX as an extension of the \LaTeX language for both authoring Linked Data vocabularies and annotating semantic documents with them. Thus, it is obviously related to other semantic extensions of \LaTeX . But, when considering \LaTeX as a text- and macro-based frontend to OWL and RDFa, it can also be compared to other ontology/vocabulary authoring and document annotation frontends, including such with graphical user interfaces.

SALT [15] also allows for annotating semantic relations in \LaTeX documents and exporting them as Linked Data. SALT is restricted to a fixed set of rhetorical and bibliographical relations, plus the metadata fields of widely used document classes like LNCS, both of which it embeds as RDF annotations in the generated PDF, whereas \LaTeX allows for (re)using arbitrary relations plus defining custom ones. The target format of \LaTeX is RDFa inside the generated OMDoc and XHTML+MathML. We have concentrated on that target, since it supports dynamic interactions via our

⁴Reuse of vocabularies is not limited by traditional restrictions of \TeX , which has a single global namespace for macros, and where no two keys passed to a command or environment may have the same name. \LaTeX groups symbols into modules; \LaTeX does the same for keys. When two symbols or keys that have the same local name relatively to their module are imported into another module M , there are facilities for giving them distinct names for usage inside M . For example, when there is already a key name, but the name property from the FOAF ontology should also be reused, we can set up a *qualified import* of the latter, e.g. as `FOAFname`.

JOBAD system. An export of the metadata relations to XMP annotations embedded in PDF should be possible with the technology employed in SALT; we leave this to future work.

SOBOLEO [6] is a lightweight graphical user interface for creating and editing vocabularies/ontologies in OWL based on Web 2.0 tagging approaches. In [5], the authors evaluate its usage along their “*Ontology Maturing Process Model*”, in which they confirm the succeeding phases “emergence of ideas”, “consolidation in communities”, “formalization”, and “axiomatization” in an ontology engineering process. Our observed phases of spotting, relating and chunking essentially correspond, as the “emergence of ideas” period did not apply (the documents were already created). Interestingly, the “consolidation in communities” phase does not only have to be thought of as a development time: We found it reified in SAMSDocs like the V-model relations. loomp is an example of a WYSIWYG editor for annotating HTML documents with terms from vocabularies, yielding RDFa [18]. GUI tools traditionally separate the task of vocabulary creation from document annotation; this also holds for SOBOLEO (responsible for the former task) and loomp (responsible for the latter). $\text{\texttt{\textit{g}\TeX+}}$, on the other hand, gives access to both tasks via the same interface: $\text{\texttt{\textit{T}\TeX}}$ macros, which are once declared, and once used – possibly even in the same source file.

7. CONCLUSION AND FUTURE WORK

We reported on a formalization case study, where we use the $\text{\texttt{\textit{g}\TeX}}$ format, a document formatting system and specification platform for semantic, mathematical vocabularies, on a document corpus from Software Engineering. To cope with the the multi-dimensional semantic structure implicit in the document collection, we extended $\text{\texttt{\textit{g}\TeX}}$ into a markup platform for semi-formal ontologies and Linked Data called $\text{\texttt{\textit{g}\TeX+}}$ (in our case semi-formal documents with RDFa-based metadata annotations).

The key observation from our case study is that if we use $\text{\texttt{\textit{g}\TeX+}}$ as a human- and document-oriented frontend for Linked Data documents, we can approach the formalization of semi-formal document collections as a process of “*document and ontology co-development*”, where (in our case pre-existing) documents are semantically preloaded with inter- and intra-document relations, whose meaning is given by (project-specific or general, reusable) metadata ontologies. As we have seen in section 3, preloading documents and developing metadata ontologies in a joint frontend format reduces formalization barriers. For instance, we often have to elaborate informal document fragments into metadata vocabularies; see the discussion about the “rd.” document state.

For practical applicability of the $\text{\texttt{\textit{g}\TeX+}}$ approach, machine support for authoring and managing $\text{\texttt{\textit{g}\TeX}}$ document collections is crucial. As a client-side counterpart to the integrated repository and Linked Data publishing solution provided by TNTBase [11], we are currently developing an integrated collection authoring environment $\text{\texttt{\textit{g}\TeXIDE}}$ for $\text{\texttt{\textit{g}\TeX}}$ on the basis of the Eclipse framework [20]. We expect that extending $\text{\texttt{\textit{g}\TeXIDE}}$ to operationalize the $\text{\texttt{\textit{g}\TeX+}}$ functionality presented in this paper will turn it into an IDE for document

collection and ontology co-development that will enable authors to cope with the complexities of dealing with large collections of semi-formalized documents. On the other hand, we expect the modular $\text{\texttt{\textit{g}\TeXIDE}}$ system to be a good basis for deploying supportive services in a flexible document collection environment.

We conjecture that the $\text{\texttt{\textit{g}\TeX+}}$ based workflow for document and ontology co-development can be extended to arbitrary Linked Data applications.

Acknowledgments. The authors gratefully acknowledge the careful work of Christoph Lüth, Holger Täubig, and Dennis Walter that went into preparing the SAMS document collection, which is the basis of this paper. Moreover, we like to thank the members of the FormalSafe project for valuable discussions.

8. REFERENCES

- [1] B. Adida and M. Birbeck. RDFa Primer. W3C Working Group Note, World Wide Web Consortium (W3C), Oct. 2008.
- [2] R. Ausbrooks, S. Buswell, D. Carlisle, G. Chavchanidze, S. Dalmas, S. Devitt, A. Diaz, S. Dooley, R. Hunter, P. Ion, M. Kohlhase, A. Lazrek, P. Libbrecht, B. Miller, R. Miner, M. Sargent, B. Smith, N. Soiffer, R. Sutor, and S. Watt. Mathematical Markup Language (MathML) version 3.0. W3C Candidate Recommendation of 15 December 2009, World Wide Web Consortium, 2009.
- [3] T. Berners-Lee, Y. Chen, L. Chilton, D. Connolly, R. Dhanaraj, J. Hollenbach, A. Lerer, and D. Sheets. Tabulator: Exploring and analyzing linked data on the semantic web. In *Proceedings of the The 3rd International Semantic Web User Interaction Workshop (SWUI06)*, Nov. 2006.
- [4] C. Bizer, T. Gauß, R. Cyganiak, and O. Hartig. Semantic web client library. <http://www4.wiwi.fu-berlin.de/bizer/ng4j/semwebclient/>, seen Feb. 2010.
- [5] S. Braun, A. Schmidt, A. Walter, and V. Zacharias. Using the ontology maturing process model for searching, managing and retrieving resources with semantic technologies. In R. Meersman and Z. Tari, editors, *OTM 2008, Part II*, number 5332 in LNCS, pages 1567–1577. Springer-Verlag, 2008.
- [6] S. Braun and V. Zacharias. SOBOLEO – a repository for living ontologies. In M. d’Aquin, A. García Castro, C. Lange, and K. Viljanen, editors, *1st Workshop on Ontology Repositories and Editors*, CEUR Workshop Proceedings, Hersonissos, Greece, May 2010.
- [7] S. Buswell, O. Caprotti, D. P. Carlisle, M. C. Dewar, M. Gaetano, and M. Kohlhase. The Open Math standard, version 2.0. Technical report, The Open Math Society, 2004.
- [8] J. Carette, L. Dixon, C. Sacerdoti Coen, and S. M. Watt, editors. *MKM/Calculus 2009 Proceedings*, number 5625 in LNAI. Springer Verlag, July 2009.
- [9] D. Carlisle. *The keyval package*. The Comprehensive $\text{\texttt{\textit{T}\TeX}}$ Archive Network, 1999. Part of the $\text{\texttt{\textit{T}\TeX}}$ distribution.

- [10] *Intelligent Computer Mathematics*, number 6167 in LNAI. Springer Verlag, 2010. in press.
- [11] C. David, M. Kohlhase, C. Lange, F. Rabe, N. Zhiltsov, and V. Zholudev. Publishing math lecture notes as linked data. In L. Aroyo, G. Antoniou, E. Hyvönen, A. ten Teije, H. Stuckenschmidt, L. Cabral, and T. Tudorache, editors, *ESWC*, number 6089 in Lecture Notes in Computer Science, pages 370–375. Springer, June 2010.
- [12] FormalSafe. <http://www.dfki.de/sks/formalsafe/>, 2008. seen March 2010.
- [13] U. Frese, D. Hausmann, C. Lüth, H. Täubig, and D. Walter. The importance of being formal. In H. Hungar, editor, *International Workshop on the Certification of Safety-Critical Software Controlled Systems SafeCert’08*, volume 238 of *Electronic Notes in Theoretical Computer Science*, pages 57–70, Sept. 2008.
- [14] J. Giceva, C. Lange, and F. Rabe. Integrating web services into active mathematical documents. In Carrette et al. [8], pages 279–293.
- [15] T. Groza, S. Handschuh, K. Möller, and S. Decker. SALT – semantically annotated L^AT_EX for scientific publications. In E. Franconi, M. Kifer, and W. May, editors, *ESWC*, number 4519 in Lecture Notes in Computer Science, pages 518–532. Springer, 2007.
- [16] O. Hartig and J. Sequeda. SQUIN – query the web of linked data. <http://squin.sourceforge.net>, seen Feb. 2010.
- [17] T. Heath et al. Linked data – connect distributed data across the web – guides and tutorials. <http://linkeddata.org/guides-and-tutorials>, seen Feb. 2010.
- [18] R. Heese, M. Luczak-Rösch, R. Oldakowski, O. Streibel, and A. Paschke. One click annotation. In T. Tudorache, G. Correndo, N. Noy, H. Alani, and M. Greaves, editors, *Proceedings of the Workshop on Collaborative Construction, Management and Linking of Structured Knowledge (CK2009)*, number 514 in CEUR Workshop Proceedings, 2009.
- [19] JOMDoc project — Java library for OMDoc documents. <http://jomdoc.omdoc.org>, 2010. seen Feb.
- [20] C. Jucovschi and M. Kohlhase. sTeXIDE: An integrated development environment for sTeX collections. In CICM10 [10]. in press.
- [21] A. Kohlhase. Overcoming Proprietary Hurdles: CPoint as Invasive Editor. In F. de Vries, G. Attwell, R. Elferink, and A. Tödt, editors, *Open Source for Education in Europe: Research and Practise*, pages 51–56, Heerlen, The Netherlands, Nov. 2005. Open Universiteit Nederland, Open Universiteit Nederland. Proceedings at <http://hdl.handle.net/1820/483>.
- [22] A. Kohlhase, M. Kohlhase, and C. Lange. Dimensions of formality: A case study for MKM in software engineering. In CICM10 [10]. <http://arxiv.org/abs/1004.5071>.
- [23] M. Kohlhase. owl2onto.cls: Marking up OWL2 Ontologies in sTeX. <https://svn.kwarc.info/repos/stex/trunk/sty/owl2onto/owl2onto.pdf>.
- [24] M. Kohlhase. RDFa metadata in L^AT_EX. <https://svn.kwarc.info/repos/stex/trunk/sty/rdfmeta/rdfmeta.pdf>.
- [25] M. Kohlhase. OMDoc – An open markup format for mathematical documents [Version 1.2]. Number 4180 in LNAI. Springer Verlag, Aug. 2006.
- [26] M. Kohlhase. Using L^AT_EX as a semantic markup format. *Mathematics in Computer Science*, 2(2):279–304, 2008.
- [27] M. Kohlhase. An open markup format for mathematical documents OMDoc [version 1.3]. Draft Specification, 2010.
- [28] M. Kohlhase, J. Giceva, C. Lange, and V. Zholudev. JOBAD – interactive mathematical documents. In B. Endres-Niggemeyer, V. Zacharias, and P. Hitzler, editors, *AI Mashup Challenge 2009, KI Conference*, Sept. 2009.
- [29] C. Lange. The OMDoc document ontology. web page at <http://kwarc.info/projects/docOnto/omdoc.html>, 2010. seen 3/2010.
- [30] C. Lange. *Semantic Web Collaboration on Semiformal Mathematical Knowledge*. PhD thesis, Jacobs University Bremen, 2010. submission expected in spring 2010.
- [31] C. Lange and M. Kohlhase. A mathematical approach to ontology authoring and documentation. In Carrette et al. [8], pages 389–404.
- [32] B. Miller. LaTeXML: A L^AT_EX to XML converter. Web Manual at <http://dlmf.nist.gov/LaTeXML/>, seen May 2010.
- [33] T. Nipkow, L. C. Paulson, and M. Wenzel. *Isabelle/HOL — A Proof Assistant for Higher-Order Logic*. Number 2283 in LNCS. Springer, 2002.
- [34] OpenLink Software. OpenLink universal integration middleware – Virtuoso product family. web page at <http://virtuoso.openlinksw.com>.
- [35] SAMS. SAMSDocs: The document collection of the SAMS project, 2009. <http://www.sams-projekt.de>.
- [36] F.-P. Servant. Linking enterprise data. In C. Bizer, T. Heath, K. Idehen, and T. Berners-Lee, editors, *Linked Data on the Web (LDOW 2008)*, number 369 in CEUR Workshop Proceedings, Apr. 2008.
- [37] Semantic Markup for LaTeX, seen July 2009. available at <http://kwarc.info/projects/stex/>.
- [38] J. Tennison et al. rdfQuery – RDF processing in your browser. <http://code.google.com/p/rdfquery/>, seen Feb. 2010.
- [39] V. Zholudev and M. Kohlhase. TNTBase: a versioned storage for XML. In *Proceedings of Balisage: The Markup Conference 2009*, Balisage Series on Markup Technologies. Mulberry Technologies, Inc., 2009. available at <http://kwarc.info/vzholudev/pubs/balisage.pdf>.
- [40] V. Zholudev, M. Kohlhase, and F. Rabe. A [insert xml format] database for [insert cool application]. In *Proceedings of XML Prague 2010*, 2010.